

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2001-357048A

(43)Date of publication of application : 26.12.2001

(51)Int.Cl.

G06F 17/30

G06F 7/24

G06F 12/00

H03M 7/30

(21)Application number : 2000-182320

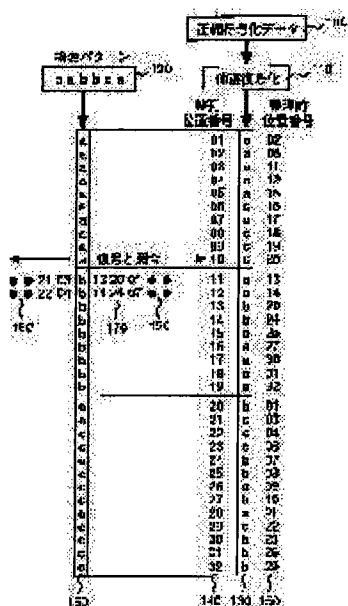
(71)Applicant : HITACHI LTD

(22)Date of filing : 13.06.2000

(72)Inventor : TONOMURA MOTONOBU

(54) METHOD FOR RETRIEVING BLOCK SORT COMPRESSED DATA AND ENCODING METHOD FOR BLOCK SORT COMPRESSION SUITABLE FOR RETRIEVAL

図 1



(57)Abstract:

PROBLEM TO BE SOLVED: To perform a high-speed retrieval by successively decoding and retrieving only required data even without decoding all encoded data concerning data compressed by a block sort compression method.

SOLUTION: Concerning data compressed by the block sort compression method, a pair of the alignment position number and the pre-alignment position number of a BW-converted column and columns aligned in a dictionary form is found. Then, collation with a retrieval data stream is performed while decoding the data on the basis of that pair. Decoding is applied only to a part required for retrieval. Besides, a pair of the alignment position number and the pre-alignment position number is directly encoded by the block sort compression

method.

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開2001-357048

(P2001-357048A)

(43) 公開日 平成13年12月26日 (2001. 12. 26)

(51) Int.Cl. ⁷	識別記号	F I	テーマコード* (参考)
G 0 6 F 17/30	2 3 0	G 0 6 F 17/30	2 3 0 A 5 B 0 7 5
	1 7 0		1 7 0 A 5 B 0 8 2
7/24		7/24	Z 5 J 0 6 4
12/00	5 1 1	12/00	5 1 1 C
H 0 3 M 7/30		H 0 3 M 7/30	Z
審査請求 未請求 請求項の数 6 O L (全 11 頁)			

(21) 出願番号 特願2000-182320(P2000-182320)

(22) 出願日 平成12年6月13日 (2000. 6. 13)

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(72) 発明者 外村 元伸

東京都国分寺市東恋ヶ窪一丁目280番地

株式会社日立製作所中央研究所内

(74) 代理人 100068504

弁理士 小川 勝男 (外2名)

Fターム(参考) 5B075 ND02 NR16 PQ02

5B082 AA11 GA02

5J064 AA02 BC29 BD04

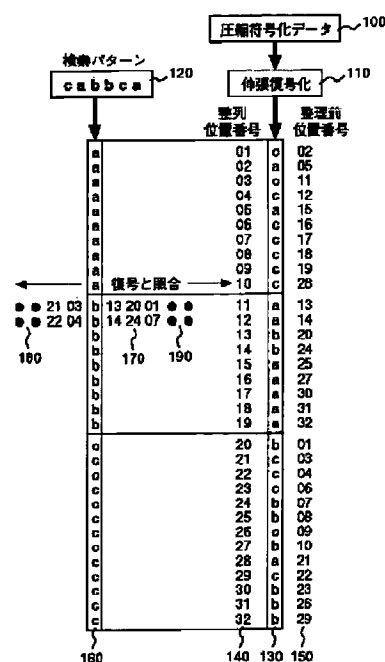
(54) 【発明の名称】 ブロックソート圧縮データの検索方法、および検索に適したブロックソート圧縮法の符号化方法

(57) 【要約】

【課題】ブロックソート圧縮法で圧縮されたデータに対して、全ての符号化されたデータを復号しなくても、逐次必要なデータだけ復号して検索することによって高速検索をおこなう。

【解決手段】ブロックソート圧縮法により圧縮されたデータに対して、BW変換した列と辞書式に整列させた列との整列位置番号と整列前位置番号の対を求める。そして、その対に基づきデータを復号しながら、検索データ列との照合をおこなう。復号は、検索のために必要な部分のみをおこなう。また、ブロックソート圧縮法で整列位置番号と整列前位置番号の対を直接に符号化しておく。

図 1



【特許請求の範囲】

【請求項1】 ブロックソート圧縮データの検索方法において、

ブロックソート圧縮法により符号化されたデータの列を第一の列、その第一の列を辞書式に整列させたデータの列を第二の列としたときに

(1) 前記第一の列のデータを前記第二の列のデータに整列させたときの整列位置番号と整列前位置番号の対を求めるステップ、

(2) 前記(1)のステップにより求められた整列位置番号と整列前位置番号の対に基づき、元のデータ列を復号するステップ、

(3) 検索データ列を前記(2)のステップにより復号されたデータ列を照合するステップとからなり、前記第一の列のデータと検索データ列を入力し、前記(1)のステップの後に、前記(2)のステップをおこない、順次復号されたデータから前記(3)のステップをおこなって元のデータ列に検索データ列が含まれるか否かを調べることを特徴とするブロックソート圧縮データの検索方法。

【請求項2】 ブロックソート圧縮法により符号化されたデータは、データの構成要素の各出現個数がわかるように、符号化されており、

前記(2)のステップにおいて、全ての元のデータ列を復号しなくても、

前記データの構成要素の各出現個数に基づいて、前記整列位置番号と整列前位置番号の対を求めて、前記(3)のステップで検索データ列との照合に必要なデータだけ復号して検索をおこなうことを特徴とする請求項1記載のブロックソート圧縮データの検索方法。

【請求項3】 ブロックソート圧縮法の符号化方法において、

巡回シフトの後に、取り出した列を符号化するときに、その取り出した列を、辞書式に整列した列に変換する場合の整列位置番号と整列前位置番号の対を直接に符号化することを特徴とするブロックソート圧縮法の符号化方法。

【請求項4】 前記(3)のステップと検索データ列と元のデータ列を照合するときに、元のデータ列の構成要素の出現個数のすくない構成要素から照合をおこなうことを特徴とする請求項1記載のブロックソート圧縮データの検索方法。

【請求項5】 前記検索データ列が一意的でない構成要素の表現がされ、

前記(3)のステップと検索データ列と元の列を照合するときに、その表現により、複数の構成要素とマッチングをすることにして検索することを特徴とする請求項1記載のブロックソート圧縮データの検索方法。

【請求項6】 前記(3)のステップにおいて、元のデータ列の中で、前記検索したデータ列を含まれる

場所の前後のデータ列をも、復号して表示することを特徴とする請求項1記載のブロックソート圧縮データの検索方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、ブロックソート圧縮データの検索方法に係り、ブロックソート圧縮法で圧縮されたデータに対して、ブロックソート圧縮法の特質を活かすことにより、全ての符号化されたデータを復号しなくても、逐次必要なデータだけ復号して検索することによって高速検索を可能にするブロックソート圧縮データの検索方法に関する。

【0002】

【従来の技術】コンピュータなどの情報機器が身近なものになり、我々がデジタルデータを取扱う機会が増すにつれて、データを符号化して圧縮して記憶し、必要に応じて、それを復号化して伸張して利用する技術が注目されている。ここで、「符号化」とは、もとのコード体系を別のコード体系に変換することであり、「復号化」とは、その逆であり、符号化したコード体系をもとのコード体系に変換し直すことと定義できる。また、「圧縮」とは、もとのデータをより少ない容量の記憶領域で格納することであり、「伸張」とは、圧縮したデータをもとのデータの容量の記憶領域を占有するようなデータにすることである。

【0003】このようなデータの圧縮／伸張技術は、パーソナルコンピュータが普及して、日常的に使われるようになってきている。中でも、1977年にZivとLempelによって提案されたLZ圧縮法は、知名度があり今日よく使われている圧縮法である。このような状況の中で、最近、LZ圧縮法の圧縮率に匹敵するブロックソート圧縮法と呼ばれる別の圧縮法がその圧縮率の高さから理論的な面で関心が集まっている(Michelle Effros, Universal Lossless Source Coding with the Burrows Wheeler Transform, IEEE Proc. of DCC'99, pp. 178-187, 1999)。

【0004】このブロックソート圧縮法は、テキスト・データ全体に対して巡回シフト(あるいは回転シフト)列を作り、すべての巡回シフト列に辞書式順序付けをおこなって配列し、そのある列を取り出して符号化するものである。例えば、この方式を発案したBurrowsとWheeler (A block-sorting lossless data compression algorithm, SRC Research Report, 124 May 1994)は、符号の列として配列の最後の列を選んでいる。

【0005】ブロックソート圧縮法は、LZ圧縮法に匹敵するような圧縮率の評価データが得られているが、まだその圧縮率の達成に関して理論的な検討段階にあり、応用面での検討があまりなされていない。

【0006】

【発明が解決しようとする課題】ところで、情報を圧縮

するばかりでなく、情報を高速に検索したいという要求がある。高速検索という目的を優先させると、どうしても検索のために冗長な情報をもたせることになるため、データ量を圧縮するというよりは、とすれば、データ量が増加することにもなる。

【0007】しかし、処理する情報が極めて大量になってくると格納スペースがなくなり、データを圧縮格納しておく必要に迫られるので、ほとんどのデータが圧縮されたままの状態で眠っているという状況になる。このような圧縮データの海から検索して必要なごく一部のデータのみを取り出すという技術が要求されることになる。すべての圧縮データを伸張復号しながら検索することは現実的な解答とはいえない。そのために、できれば圧縮コード列のまま伸張しないで検索する方法を見出したいところである。

【0008】しかしながら実際には、既存のLZ圧縮法によるデータに対して、圧縮コード列に対して検索パターンと比較照合する場合、圧縮コード内容の前後の一部にまたがって照合したり、検索パターンの圧縮コード列が一意に定まらず、何通りも存在することがあるため、圧縮コード列のままでは直接検索することができないという問題がある。ブロックソート圧縮法は、評価段階の圧縮技術と言って良く、これまで、ブロックソート圧縮データに対する検索法について十分考察されてこなかった。

【0009】本発明は、上記問題点を解決するためになされたもので、その目的は、ブロックソート圧縮法で圧縮されたデータに対してブロックソート圧縮法の特質を活かすことにより、全ての符号化されたデータを復号しなくても、逐次必要なデータだけ復号して検索することによって高速検索を可能にするブロックソート圧縮データ検索法を提供することにある。

【0010】また、その検索に適したブロックソート圧縮法の符号化方法を提供することにある。

【0011】

【課題を解決するための手段】ブロックソート圧縮法は、テキスト・データ全体に対して巡回シフト列を作り、すべての巡回シフト列に辞書式順序付けをおこなって配列化する。それで、検索したいパターンが複数箇所にある場合、巡回シフト列の配列の中で、検索パターンは配列のある行の先頭列から始まり、しかも複数箇所の検索パターンが配列の連続した行において固まって出現するという特徴がある。また、ブロックソート圧縮データの復号原理において、伸張復号した最後の列の文字列の位置を辞書式順序で並べ替えて整列させるが、そのとき、整列位置番号と整列前の整列位置番号の対を作り、元のテキストの整列位置番号を指定して、これらの対を順次たどりながら、テキストの先頭から順番に復号できる。

【0012】したがって、本発明では、ブロックソート

圧縮データがもつこのような性質を利用して効率良い検索手段を提供する。すなわち、まず検索パターンの先頭の文字と2番目の文字の対に対して、整列位置番号と整列前位置番号の対を当てはめて対応させる。これらに当てはめられる対は、辞書式順序で整列させているために、固まって出現するので候補を絞り込める。続いて検索パターンの2番目の文字と3番目の文字の対に対して、前段で絞り込まれたものに対して整列位置番号と整列前の整列位置番号の対を当てはめて対応させるという手順を同様に繰り返していく。すると、検索パターンの長さをnとすると、検索パターンのn-1番目の文字とn番目の文字の対に対して整列位置番号と整列前位置番号の対を当てはめて絞り込み対応させた段階で一連の手順は終了する。

【0013】結果として、複数箇所で見出された検索パターンだけが残し、同時に元のデータ列に含まれている検索パターンが見出されることになる。

【0014】元のテキスト列の文字の出現個数がわかっているときには、整列位置番号と整列前位置番号の対を順次求めることができるため、照合にあたり全ての元のテキスト列を復号する必要はなく、検索パターンと照合する必要のあるところだけ復号して照合すれば良い。また、いわゆるあいまい検索も、マッチングの可能性のあるところを上記手順で復号することによりおこなうことができる。

【0015】検索パターンを照合するときには、元のテキスト列で出てくる文字が一番少ない文字からおこなえば絞込みが速くなり、効率的に検索できる。

【0016】ブロックソート圧縮法では、符号化操作は2段階でおこなわれる。第1段階の符号化では、通常の発想では、連続して出てくる文字の長さに着目して符号化するものである。しかしながら、上記検索方法では、第1段階の復号化のステップと整列位置番号と整列前位置番号の対を求める手順が無関係に存在しており、効率的には改善すべきものがある。

【0017】そこで、ブロックソート圧縮法において、配列の最後の列の文字列を圧縮符号化するのではなく、整列位置番号と整列前位置番号の対を直接に圧縮符号化することで復号と検索の効率をさらに上げることとする。整列位置番号と整列前位置番号の対は、配列の最後の列の文字列に1対1で対応しているため、ほぼ同程度の圧縮率の達成が期待できる。ブロックソート圧縮法の符号化方法によって、検索向きの圧縮法を提供することができる。

【0018】

【発明の実施の形態】以下、本発明に係る各実施形態を、図1ないし図11を用いて説明する。

【ブロックソート圧縮法】先ず、本発明のブロックソート圧縮データの検索法を説明する前に、図2および図3を用いて前提となるブロックソート圧縮法について説明

する。図2は、ブロックソート圧縮法の圧縮符号化過程を説明するための概略図である。図3は、ブロックソート圧縮法を具体的なデータにより説明するための図である。

【0019】以下、本発明の実施形態では、一貫して、元のテキストが'cabccabcccabbbcabccabccabbbcaab'の32個の文字からなるテキスト200を圧縮した場合を例に採り説明していくことにする。

【0020】ブロックソート圧縮法により圧縮して、符号化するときのアルゴリズムの概略は、以下のようになる。

【0021】(圧縮-ステップ1) 先ず、元のテキスト200の巡回シフトをおこなって、全ての巡回シフト列210を求める。巡回シフトとは、元の列を一字単位で右または左に、回転させてシフトさせることであり、図2の例では、テキスト200を左側にシフトさせて、先頭からはみ出た文字'c'が、最後尾にきている。

【0022】この元のテキスト200の例では、32個の文字からなっているため、32個の巡回シフト列210ができることになる。

【0023】(圧縮-ステップ2) (圧縮-ステップ1)で生成された巡回シフト列を辞書式順序で整列させて、配列220を作る。

【0024】(圧縮-ステップ3) 配列220の最後の列130を取り出して、これを圧縮符号化する。このように元のテキスト200からこのような手順を経て、最後の列130の列に変換することを、発案者の名前をとり、Burrows-Wheeler変換、略してBW変換と言っている。実際には、取ってくるのは配列220のどの列でも良いが、前記論文には、最後の列をとっている。

【0025】また、この配列220の中で、元のテキストの位置番号230である'25'についても圧縮しておく。

【0026】元のテキスト200とBW変換した列は、長さは同じであるが同じ文字が続く傾向があることが知られているので、例えば、文字列の連続長を符号化することにより高い圧縮率が得られる。なお、BW変換した列の符号化の仕方は、外にもいろいろな方法が考えられ、必ずしも上記の方法にこだわる必要はない。

【0027】このように、ブロックソート圧縮法は、辞書式順序で整列させた配列220に基づいて、符号化のためのデータを得るものなので、元のテキスト200を圧縮するよりも効率の良い圧縮がおこなえるのでないかとして注目されているものである。

【0028】次に、上記手順により、圧縮して符号化されたデータを復号して伸張する場合の手順を、図4および図5を用いて説明する。図4は、ブロックソート圧縮法の伸張復号化過程を説明するための概略図である。図

5は、ブロックソート圧縮法を伸張復号化過程を具体的なデータにより説明するための図である。

【0029】先ず、具体的な手順を説明する前に、図3に示されている整列位置番号と整列前位置番号について説明する。この整列位置番号と整列前位置番号は、ブロックソート圧縮法のアルゴリズムを理解する上において、非常に重要な概念である。

【0030】整列位置番号は、巡回シフト列を辞書式順序に整列させた配列220の位置そのものである。

【0031】整列前位置番号とは、BW変換した最後の列130を配列の最初の列に一致させるように、整列させたときに、整列させた文字が、整列する前にはどの整列位置番号の位置にあったかを示す位置番号である。

【0032】具体的には、BW変換した最後の列130は、'caccaccccccaabbaaaaaabcccbcbcbacbb'である。

【0033】これを整列させるために、先ず'a'の文字が来る。これは、整列前には、2番目にあったものだから、整列前位置番号は、02となる。次に来るのも、'a'の文字である。次の'a'の文字が見出されるのは、BW変換した最後の列130で5番目なので、整列前位置番号は、05となる。

【0034】同様にして、'a'の文字が並び、'b'の文字が最初に来るときの整列前位置番号は、13になり、'c'の文字が最初に来るときの整列前位置番号は、01になる。このようにして、整列した列130が得られる。この対応の原理は、図4によって示されている。

【0035】ここで、記号の約束をする。整列位置番号140と整列前位置番号150の対を、(整列位置番号、整列前位置番号)と書くことにする。例えば、'a'が最初に来るところは、(01, 02)、'b'が最初に来るところは、(11, 13)、'c'が最初に来るところは、(20, 01)である。

【0036】次に、上記手順により圧縮符号化したときのデータを、伸張復号化するときのアルゴリズムの概略は、以下のようになる。

【0037】(伸張-ステップ1) (圧縮-ステップ3)で符号化した元のテキスト位置230と、BW変換した最後の列130を伸張復号化する。これを仮に第1段階の伸張復号化ということにする。この第1段階の伸張復号化は、もちろん、(圧縮-ステップ3)で符号化したアルゴリズムに基づくものである。

【0038】これにより、元のテキスト位置230である'25'と、BW変換した最後の列130、'caccaccccccaabbaaaaaabcccbcbcbacbb'が得られたものとする。

【0039】(伸張-ステップ2) (伸張-ステップ1)で得られたBW変換した最後の列130を、辞書式で整列させる。そのとき、上記手順で得られるような(整列位置番号、整列前位置番号)の対も記憶してお

く。

【0040】この例では、図3に示されるように(01, 02)、(02, 05)、(03, 11)、…、(32, 29)という対を得ることができる。

【0041】(伸張-ステップ3)元のテキスト位置230、BW変換した最後の列130、整列させた列160および(整列位置番号、整列前位置番号)を基にして元のテキスト200を復号する。これが第2段階の伸張復号化である。

【0042】この例で示すと、以下のようになる。

【0043】まず、元のテキスト位置230が'25'であるから、整列させた列160(図3の最初の列)の25番目を見て、最初の文字、'c'が復号される。この'c'は、整列前には、8番目にあったことが、(25, 08)の対を見ることにより分る(図5①)。次に、整列させた列160の8番目は、'a'である。もともと、この配列は巡回シフト列なのだから、この'a'は、最初の文字'c'の次に来るものである。したがって、2番目の文字'a'が復号される(図5②)。

【0044】同様に、(08, 18)の対を見て、整列させた列160の18番目に、'b'が来ているので、3番目の文字'b'が復号される。

【0045】このような元のテキスト位置230から、'25'から(25, 08)、(08, 18)、(18, 31)、(31, 26)、…というように連鎖的にたどっていくと、順々に、'c a b c...'というように元のテキスト200が復号化されて得られることになる。

【0046】このようにブロックソート圧縮法は、巡回シフト列の性質を巧妙に利用して圧縮伸張をおこなうものである。

【0047】〔ブロックソート圧縮データの検索方法の基本原則〕次に、図6を用いてブロックソート圧縮法により圧縮されたデータ(以下単に、「ブロックソート圧縮データ」という)に対して、特定のパターンを検索するための基本原則について説明する。図6は、本発明に係るブロックソート圧縮データの検索方法の検索過程を説明するための図である。

【0048】本実施形態では、検索パターン120として、'c a b b c a'を取り上げる。この検索パターンは、元のテキスト200には、2箇所に見出される。

【0049】ここで、記号として検索パターンのi番目の文字を、P[i]で表すことにする。この例では、図6にも示されているようにP[1]='c'、P[2]='a'などである。

【0050】アルゴリズムとしては、まず、検索パターン120の最初の文字P[1]が、整列した列160(最初の列)のどこに見い出されかをサーチする。この整列した列160は、辞書式に整列しているため、その文字が最初に出現した場所から連続して出てくる個数だけ見れば良く、サーチは極めて容易である。元のテキス

ト200から直接サーチする場合には、先頭の文字から始めて、順番に照合しなければならず、サーチが元のテキスト200の長さ分だけおこなわなければならない。したがって、それと比較すると、最初から整列した列からサーチするためこの検索は極めて効率的であるといえることができる。

【0051】図6の表では、P[1]='c'の下に、20~32までの数字が並んでいるが、これが、整列した列160の整列位置番号である。実際に図3では、20~32に'c'が来ていることが確認できる。

【0052】次に、P[2]='a'を検索する。

【0053】これは、各々P[1]で見つかった整列位置番号から整列前位置番号の対を求め、それにより検索する。すなわち、整列位置番号'20'から整列前位置番号'01'を求め、ブロックソート圧縮法の復号の原理により、'a'が復元されて、2文字目も一致することが分る。このように2文字目が'a'のパターンを調べると、図6の表の2列目から分るように、(整列位置番号、整列前位置番号)の組は、(20, 01)、(21, 03)、(22, 04)、(23, 06)、(24, 07)、(25, 08)、(26, 09)、(27, 10)となる。次の整列位置番号'28'では、整列前位置番号'18'であり、復号される文字は、'c'となるため検索パターンと一致しないことが分る。そして、これ以降では、検索パターンと一致するパターンは原理上見出されない。というのも、配列220は、もともと辞書式に整列しているからである。

【0054】同様に、次のP[3]='b'を検索する。P[2]の整列位置番号から見出される候補は、(03, 11)、(04, 12)、(06, 16)、(07, 17)、(08, 18)、(09, 19)である。

【0055】このようにして、P[1]~P[6]まで、一致するのは、図6に示されるように(21, 03)、(03, 11)、(11, 13)、(13, 20)、(20, 01)の行610と(22, 04)、(04, 12)、(12, 14)、(14, 24)、(24, 07)の行620の二箇所で見出されることが分る。すなわち、この場所で検索する6個の文字が見つかったことを意味する。

【0056】この検索法によれば、ブロックソート圧縮法の整列した列160と(整列位置番号、整列前位置番号)を基にして、先頭位置P[1]から順番に調べて行けば良く、しかも、各々の探索枝に対して、必ず検索パターンの長さだけ探せば良いので、元のテキスト列200をサーチする場合に比べて極めて能率的な検索をおこなうことができる。

【0057】〔該当箇所の前後の表示〕ところで、元のテキスト列200から検索パターン120を検索するとき、その該当する箇所の前後の文字列を表示したいことが実用上よくある。

10

20

30

40

50

【0058】この場合においても、本発明のブロックソート圧縮データの検索方法によれば、検索したときと、同様の手順により、該当箇所の前後の文字列を復号して表示することができる。

【0059】例えば、上記の例で、図6の行610の箇所の前の文字列を表示したいとする。この場合には、整列前位置番号'21'のときの整列位置番号を求めれば、P[1]の前の文字を求めることができる。すなわち、(整列位置番号, 整列前位置番号) = (x, 21)にあたるxを求めれば良い。xは28となり、整列した列160の28番目の文字が'c'であることにより、求める文字は、'c'であることが分る。その前の文字も同様にして、(x, 28)から、xは10となり、求める文字は、'a'であることが分る。

【0060】逆に、図6の行610の箇所の後の文字列を表示したい場合には、一番最後の対(20, 01)に注目して、整列位置番号'01'のときの整列前位置番号を求めれば良い。すなわち、(整列位置番号, 整列前位置番号) = (01, y)にあたるyを求める。yは02であり、整列した列160の2番目の文字が'a'であることから、このP[6]のすぐ後の文字は、'a'であることが分る。同様に、(02, y)のyを求めると、yは05となり同様に次の文字も'a'であることができる。

【0061】このように元のテキスト200で検索パターン120がある所の前後の文字列は(整列位置番号, 整列前位置番号)の連鎖をたどっていくことにより、自然と復号でき、これをCRTやプリンタなどの出力装置に表示させることができる。

【0062】[あいまい検索への応用]次に、図7を用いて本発明に係るブロックソート圧縮データの検索方法が、あいまい検索へも応用できることを説明する。図7は、本発明に係るブロックソート圧縮データの検索方法であいまい検索をおこなったときの検索過程を説明するための図である。

【0063】テキスト検索において、いわゆる「あいまい検索」をおこないたいことが良くある。あいまい検索とは、例えば、単語の一部のみを指定し、その他の部分はなにが来ても良いとして一致するパターンを検索することである。例えば、'*'の文字がドントケアの文字(ワイルドカードとも言う)を表すとして、検索パターンに'*'を指定したときには、全ての文字とマッチングするものと約束する。

【0064】この例で、例えば、あいまい検索として、検索パターンとして'ca**ac'が指定されたものとする。すなわち、P[3]=P[4]='*'で他の部分は、上の例と同様である。

【0065】この場合には、先ず、既に述べた本発明の検索方法により、P[1]P[2]='ca'の部分で一致する箇所を検索する。一致する部分は、(整列位

置番号, 整列前位置番号)の表現で表すと、図7に示される通り、(20, 01)、(21, 03)、(22, 04)、(23, 06)、(24, 07)、(25, 08)、(26, 09)、(27, 10)の8個である。

【0066】次のP[3]P[4]='**'の部分は、全ての文字とマッチングするので、そのまま、(整列位置番号, 整列前位置番号)の連鎖をたどることになる。この過程では、候補は減らずに推移する。そして、この候補の中から後半のパターンP[5]P[6]='ca'にマッチングするもののみを追跡し、最終的な絞り込みをおこなう。すなわち、次のP[5]のところでマッチングするのは、図のように5箇所であり、P[6]のところでさらに絞り込まれ、このあいまい検索の解としては図7に示されるように4箇所の場所が求まることになる。

【0067】[ブロックソート圧縮データの検索方法の効率化—その一—]上で、本発明に係るブロックソート圧縮データの検索方法の原理について述べたが、ここでは、図8を用いてさらに本発明の検索方法を効率的におこなうための工夫について説明する。図8は、元のテキスト列の出てくる個数に着目して復号と検索をおこなうための過程を説明するための図である。

【0068】上記ブロックソート圧縮データの検索方法の原理では、ブロックソート圧縮法における第1段階の伸張復号化をおこなって、復号化されたBW変換された列130(最後の列)に基づいて検索をおこなうものとして説明した。

【0069】次に説明する本発明の検索方法は、必ずしも完全にBW変換された列130を復号しなくても、検索をおこなえるものである。したがって、一層効率的な検索がおこなえることが期待される。

【0070】この検索をおこなえるための条件は、元のテキスト200の文字の出現個数がわかるように符号化されていることである。この例では、'a'が10個、'b'が9個、'c'が13個である。これから説明する検索方法の効率化のポイントは、文字の出現個数が分っているときには、BW変換した列130の先頭から順に処理していく毎に、整列位置番号と整列前位置番号の対を求めることができるため、先頭から順に復号して検索パターンとマッチング処理をおこなうことができる点にある。

【0071】具体的に手順を追っていくと以下の通りである。先ず、BW変換された列130の最初の文字は、'c'である。これは、'c'の1番目であり、しかも、予め文字の出現個数が分っているので、'c'の整列位置番号は、10+9+1となる。すなわち、(整列位置番号, 整列前位置番号) = (20, 01)が求まる。

【0072】図8では、文字毎に整列前位置番号を並べたものであり、'a'の順番1のセルが整列位置番号が1、'b'の順番1のセルが整列位置番号が11、'c'の

順番1のセルが整列位置番号が20に該当することを示している。

【0073】そして、次の文字'a'は、'a'の1番目であり、整列位置番号は、1である。すなわち、(整列位置番号, 整列前位置番号) = (01, 02)である。同様に、3番目の文字は、'c'については、整列位置番号は、 $10 + 9 + 2 = 22$ であり、(整列位置番号, 整列前位置番号) = (22, 03)である。これは、'c'の順番2のセルの整列前位置番号が3であることに対応する。このように図8は、'a'、'b'、'c'の各文字が出てくるたびに、整列前位置番号を対応する行のセルに入れていけば、自動的に整列位置番号を求めることができることを示している。

【0074】さて、検索パターン120は、'c a b b c a'であった。

【0075】ここで、文字'b'の最初のものが出てくる場所まで、整列操作がおこなわれたとする。図8からすぐ見て取れるように、文字'b'は、最初から13番目に出てきて、すなわち、整列前位置番号は、13であり、整列位置番号は、 $10 + 1 = 11$ である。

【0076】ここまで、(整列位置番号, 整列前位置番号)の対で、検索パターン120とマッチングするものは、(21, 03) (03, 11) (11, 13)と、(22, 04) (04, 12)のシーケンスで、'c a b b'と'c a b'までが照合できる。

【0077】このようにして、BW変換した列130の24番目の文字'b'まで、整列操作がおこなわれた段階では、(21, 03) (03, 11) (11, 13) (13, 20) (20, 01)と(22, 04) (04, 12) (12, 14) (14, 24) (24, 07)で、検索パターン120の'c a b b c a'が照合できる。

【0078】そして、これ以降では、もはや検索パターン120'c a b b c a'は、出現しないことが分る。これは、(整列位置番号, 整列前位置番号) = (16, x x)にあたるx xには、文字'b'にあたる整列前位置番号11~19が来ることはない。これは、24番目まで調べているので、他の場所に既に整列前位置番号11~19が使われていることが判明しているからである。したがって、これ以降には文字'b'が来ることはなく、これ以降の照合操作はおこなう必要のないことがわかる。

【0079】通常のテキストとの照合処理により、検索処理をおこなう場合には、最後の文字まで、照合しないと検索パターンを検出できないのと比較して、本発明のブロックソート圧縮データの検索方法の有利な点である。ただし、最悪の場合には最後まで整列操作をおこなわなければならない場合も生じうる。

【0080】〔ブロックソート圧縮データの検索方法の効率化—その二〕次に、図1を用いて本発明の検索方法を効率的におこなうための他の工夫について説明する。

図1は、本発明のブロックソート圧縮データの検索方法の概略手順を説明するための図である。

【0081】これまで説明してきた本発明のブロックソート圧縮データ検索では、検索パターン120の先頭から照合処理をおこなってきた。しかしながら、検索パターン120の先頭文字P[1]のテキスト200での出現個数が多いと、最初の照合処理が多くなり絞り込むまでの操作も多くなる。したがって、これを避けるためには、検索パターン120の文字列の中で、テキスト200の中で出現個数の少ない文字を選び出し、その位置から照合処理を開始して、絞り込んだ後で、開始位置の前の文字を逆に遡って照合していくと効率的になる。

【0082】また、検索パターン120の複数箇所の出現位置を同時に検出するよりも、先ず1箇所目を見つけた方が、2箇所目以降の絞り込みがはやくなる。

【0083】検索パターン120'c a b b c a'の例では、'a'、'b'、'c'の三種類の文字が出てくるわけであるが、元のテキスト200の中で、文字'b'が9個で一番少ないので検索パターン120の3番目の文字'b'に着目して検索を始める。

【0084】図1に示されるように、前向き照合では、(11, 13) (13, 20) (20, 01)と照合されていき、後向き照合では、(21, 03) (03, 11)というシーケンスで照合されていくことになる。このように検索パターンの任意の文字を選び出して、照合処理をおこなっていけるというのも、整列位置番号と整列前位置番号によって、前向きでも後向きでもまったく対称的に、復号ができるというブロックソート圧縮法の特徴によるものである。

【0085】〔ブロックソート圧縮法の修正した圧縮符号化〕これまで、ブロックソート圧縮法を前提とする検索方法を述べてきた。この検索方法では、ブロックソート圧縮法の第1段階の伸張復号化をおこない、しかる後に(整列位置番号, 整列前位置番号)により、第2段階の伸張復号化をおこなって検索パターンとの照合をおこなうものであった。第1段階の符号化の例としては、例えば、文字列の連続する長さで符号化する方法をあげた。

【0086】ここでは、第1段階の符号化の段階で、(整列位置番号, 整列前位置番号)の対を直接符号化することにより、2段階の符号化、復号化を1段階で済ませて、さらに効率的なブロックソート圧縮データに対する検索をおこなうアイデアについて説明する。

【0087】以下、図9および図10を用いてこれまでの例により説明することにする。図9は、本発明のブロックソート圧縮法の修正した圧縮符号化を説明するための図である。図10は、圧縮符号化したデータを模式的に示した図である。

【0088】ブロックソート圧縮法は、(整列位置番号, 整列前位置番号)の対を利用して、第2段階の伸張

復号化をおこなうものであった。したがって、(整列位置番号、整列前位置番号)を直接に符号化すれば、復号時のこの対応付けを省略できるというのが基本的な発想である。

【0089】図9は、'a'のテーブル410、'b'のテーブル420、'c'のテーブル430毎に整列位置番号340と整列前位置番号350を対応してあげたものである。

【0090】整列位置番号340も整列前位置番号350も0から始めている。これは、できるだけ符号化時に記憶容量を減らそうという技術的な工夫である。

【0091】また、BW変換した列160は、おなじ文字が続く傾向があるため、整列前位置番号が同一の連番が続くことが期待される。したがって、整列前位置番号350は、各テーブルの相対位置であらわせれば、圧縮率が高くなることが予想される。したがって、整列前位置番号350は、テーブルインデックス440と一緒に、各テーブルの相対番号として表現することにする。

【0092】'a'のテーブル410の最初のエントリでは、整列位置番号340が00であり、整列前位置番号350が01であり、テーブルインデックス440が'a'になっている。これは、図3の(整列位置番号、整列前位置番号)としては、(01, 02)に該当する。また、'a'のテーブル410の3番目のエントリは、整列位置番号340が02であり、整列前位置番号350が00であり、テーブルインデックス440が'b'になっている。図8で示したように、'b'のテーブルの最初の位置は、11番目を表しているので、図3の(整列位置番号、整列前位置番号)としては、(03, 11)である。

【0093】これを実際に符号化するときには、整列前位置番号350と整列位置番号340の差を求めて相対的に符号化することにする。そして、これをテーブルインデックスと共に復号可能なように符号化する。

【0094】このようにして符号化されたデータを符号化の仕方が良くわかるような書き方をすると、図10のようになる。これは、テーブルインデックスと相対位置を符号化し、さらに、連続する文字の現れる場合の表記を工夫したものである。この表記で*i* + *j*と書いてあるのは、*i*が*j*個連続していることを示している。

【0095】この図10で、a(1, 3)は、テーブルインデックスが'a'で、整列前位置番号350と整列位置番号340の差360が、1と3、次のb(-2+, 0+4)は、テーブルインデックスが'b'であり、差360が-2, -2と続き、0が4つ続くことを意味している。

【0096】〔ブロックソート圧縮データの検索方法のアルゴリズム〕最後に、図1を用いてこれまで説明してきたことを基にしてブロックソート圧縮データの検索方法のアルゴリズムを整理して述べることにする。図1

は、本発明のブロックソート圧縮データの検索方法の概略手順を説明するための図である。

【0097】元のテキスト200がブロックソート圧縮法によって圧縮符号化されたデータが記憶媒体に記憶されているとする。

【0098】そして、検索をおこなう検索パターン120が指定されているとする。

【0099】本発明の検索方法は、圧縮符号化データ100を伸張復号化しながら、検索パターン120との照合処理を可能にするものであった。

【0100】検索は、任意の文字からおこなえるのであるが、元のテキストの出現個数の一番少ない文字、この例では'b'からおこなうのが効率的である。復号化するに際しては、(整列位置番号140、整列前位置番号150)の対を順次連鎖的にたどって、検索パターンとマッチングするテキストの部分列を絞り込みながら、テキストの前後の双方向にわたって復号していく。第1段階の符号化で、出現する文字の個数がわかるようにしておくか、(整列位置番号140、整列前位置番号150)の対自体を符号化するようにすれば、圧縮符号化データ100の全てを元のテキスト200として復号しなくても順次、検索がおこなえることはこれまで述べた通りである。

【0101】検索が成功して、該当箇所が見つかったら、必要に応じてマッチング箇所の前部と後部がどのように並びになっているかをユーザに提示してやる。

【0102】

【発明の効果】本発明のブロックソート圧縮データの検索方法は、検索対象の文字列パターンがテキスト中の数箇所に出現している場合でも、すべてを対象に先頭から同時に照合することができる。しかも、検索パターンの長さ分の照合を終了したところで、すべての照合可能な位置が検出される。そのため、効率的な高速検索ができるデータの高効率圧縮法である。また、検索パターンの前後の文字列の復号が直接おこなえるので、画面に検索検出位置の前後の文字列を同時表示することができるので応用面でも便利である。また、ブロックソート圧縮法の整列位置番号と整列前位置番号の対を直接符号化しておけば、この検索方法に有効に利用できる。

【0103】このように本発明によれば、ブロックソート圧縮法で圧縮されたデータに対してブロックソート圧縮法の特質を活かすことにより、全ての符号化されたデータを復号しなくても、逐次必要なデータだけ復号して検索することによって高速検索を可能にするブロックソート圧縮データ検索法を提供することができる。

【0104】また、その検索に適したブロックソート圧縮法の符号化方法を提供することができる。

【図面の簡単な説明】

【図1】本発明のブロックソート圧縮データの検索方法の概略手順を説明するための図である。

【図2】ブロックソート圧縮法の圧縮符号化過程を説明するための概略図である。

【図3】ブロックソート圧縮法を具体的なデータにより説明するための図である。

【図4】ブロックソート圧縮法の伸張復号化過程を説明するための概略図である。

【図5】ブロックソート圧縮法を伸張復号化過程を具体的なデータにより説明するための図である。

【図6】本発明に係るブロックソート圧縮データの検索方法の検索過程を説明するための図である。

【図7】本発明に係るブロックソート圧縮データの検索方法であいま検索をおこなったときの検索過程を説明するための図である。

【図8】元のテキスト列の出てくる個数に着目して復号と検索をおこなうための過程を説明するための図である。

【図9】本発明のブロックソート圧縮法の修正した圧縮符号化を説明するための図である。

*

*【図10】圧縮符号化したデータを模式的に示した図である。

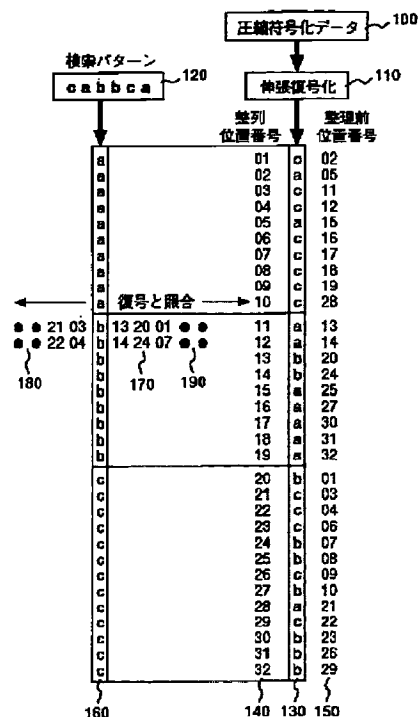
【符号の説明】

100…圧縮符号化データ、110…第1段階の伸張復号化、120…検索パターン、130…BW変換された(配列の最後列の)列、140…整列位置番号、150…整列前位置番号、160…整列させた(配列の先頭列の)列、170…検索された整列番号、180…検索パターンの前部の復号化、190…検索パターンの後部の復号化、200…元のテキスト、210…巡回シフト列、220…巡回シフト列の整列配列、230…元のテキストの位置番号、500…第2段階の復号化、340…整列位置番号(00から始まる)、350…整列前位置番号(00から始まる)、360…整列前位置番号と整列位置番号の差、410…'a'テーブル、420…'b'テーブル、430…'c'テーブル、440…テーブルインデックス。

【図1】

図 1

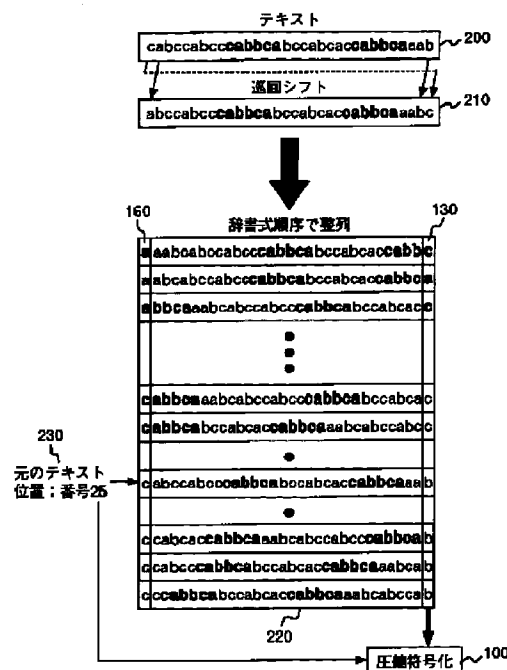
Fig 1



【図2】

図 2

Fig 2



【図8】

図 8

Fig 8

番号	a	b	c
1	02	13	01
2	05	14	08
3	11	20	04
4	12	24	06
5	15	26	07
6	18	27	08
7	17	30	09
8	18	31	10
9	19	32	21
10	28		22
11			23
12			26
13			29

Fig 9

【図9】

図 9

